

# Frankenstein: Learning Deep Face Representations using Small Data

Guosheng Hu, *Member, IEEE*, Xiaojiang Peng, *Member, IEEE*, Yongxin Yang, *Student Member, IEEE*, Timothy Hospedales, *Member, IEEE*, and Jakob Verbeek, *Member, IEEE*

**Abstract**—Deep convolutional neural networks have recently proven extremely effective for difficult face recognition problems in uncontrolled settings. To train such networks, very large training sets are needed with millions of labeled images. For some applications, such as near-infrared (NIR) face recognition, such large training datasets are, however, not publicly available and very difficult to collect. In this work, we propose a method to generate very large training datasets of synthetic images by compositing real face images in a given dataset. We show that this method enables to learn models from as few as 10,000 training images, which perform on par with models trained from 500,000 images. Using our approach we also improve the state-of-the-art results on the CASIA NIR-VIS2.0 heterogeneous face recognition dataset.

**Index Terms**—face recognition, deep learning, small training data

## I. INTRODUCTION

In recent years, deep learning methods, and in particular convolutional neural networks (CNNs), have achieved considerable success in a range of computer vision applications including object recognition [20], object detection [9], semantic segmentation [31], action recognition [39], and face recognition [35]. The recent success of CNNs stems from the following facts: (i) big annotated training datasets are currently available for a variety of recognition problems to learn rich models with millions of free parameters; (ii) massively parallel GPU implementations greatly improve the training efficiency of CNNs; and (iii) new effective CNN architectures are being proposed, such as the Very Deep VGG Network [40], Google Inception Network [49] and Deep Residual Networks [11].

Good features are essential for object recognition, including face recognition. Conventional features include linear functions of the raw pixel values, including Eigenface (Principal Component Analysis) [51], Fisherface (Linear Discriminant Analysis) [3], and Laplacianface (Locality Preserving Projection) [12]. Such linear features were later replaced by hand-crafted local non-linear features, such as Local Binary Patterns [1], Local Phase Quantisation (LPQ) [2], and Fisher vectors computed over dense SIFT descriptors [38]. Note that the latter is an example of a feature that also involves unsupervised learning. These traditional features achieve promising face recognition rates in constrained environments, as represented

in the CMU PIE dataset [37]. However, using these features face recognition performance may degrade dramatically in uncontrolled environments, as represented in the Labeled Faces in the Wild (LFW) benchmark [14]. To improve the performance in such challenging settings, metric learning can be used, see [4], [10], [52]. Metric learning methods learn a (linear) transformation of the features that pulls the objects that have the same label closer together, while pushing the objects that have different labels apart.

Although hand-crafted features and metric learning achieve promising performance for uncontrolled face recognition, it remains cumbersome to improve the design of hand-crafted local features (such as SIFT [23]) and their aggregation mechanisms (such as Fisher vectors [34]). This is because the experimental evaluation results of the features cannot be automatically fed back to improve the robustness to nuisance factors such as pose, illumination and expression. The major advantage of CNNs is that all processing layers, starting from the raw pixel-level input, have configurable parameters that can be learned from data. This obviates the need for manual feature design, and replaces it with supervised data-driven feature learning. Learning the large number of parameters in CNN models (millions of parameters are rather a rule than an exception) requires very large training datasets. For example, the CNNs, which achieve state-of-the-art performance on the LFW benchmark, are trained using datasets with millions of labeled faces: Facebook’s DeepFace [50] and Google’s FaceNet [35] were trained using 4 million and 200 million training samples, respectively.

For some recognition problems large supervised training datasets can be collected relatively easily. For example the CASIA Webface dataset [55] of 500,000 face images was collected semi-automatically from IMDb. However, in many other cases collecting large datasets may be costly, and possibly problematic due to privacy regulation. For example, thermal infrared imaging is ideal for low-light nighttime and covert face recognition applications [19], but it is not possible to collect millions of labeled training images from the internet for the thermal infrared domain. The lack of large training datasets is an important bottleneck that prevents the use of deep learning methods in such cases, as the models will overfit dramatically when using small training datasets [13].

To address this issue, the use of big synthetic training datasets has been explored by a number of authors [15], [28], [32]. There are two important advantages of using synthetic data (i) one can generate as many training samples as desired, and (ii) it allows explicit control over the nuisance factors.

G. Hu, X. Peng and J. Verbeek are with THOTH team, INRIA Grenoble Rhone-Alpes, France. email: {guosheng.hu, xiaojiang.peng, jakob.verbeek}@inria.fr. G. Hu is also with Anyvision group.

Y. Yang and T. Hospedales are with Electronic Engineering and Computer Science, Queen Mary University of London, UK. email: {yongxin.yang, t.hospedales}@qmul.ac.uk

For instance, we can synthesize face images of all desired viewpoints, whereas data collected from the internet might be mostly limited to near frontal views. Data synthesis has successfully been applied to diverse recognition problems, including text recognition [15], scene understanding [28], and object detection [32]. Two very recent works [58], [8], [25] proposed 3D-aided face synthesis technique for facial landmark detection and face recognition in the wild respectively.

Data augmentation is another technique that is commonly used to reduce the data scarcity problem, see [30], [40]. This is similar to data synthesis, but more limited in that existing training images are transformed without affecting the semantic class label, e.g. by applying cropping, rotation, scaling, etc.

The main contribution of this paper is a solution for training deep CNNs using very small training data. To achieve this, we propose a data synthesis technique to expand very limited face datasets to very large ones that are suitable to train powerful deep CNNs. Specifically, we synthesize images of a ‘virtual’ subject  $c$  by compositing automatically detected face parts (eyes, nose, mouth) of two existing subjects  $a$  and  $b$  in the dataset in a fixed pattern. Images for the new subject are generated by composing a nose from an image of subject  $a$  with a mouth of an image of subject  $b$ . This is motivated by the observation that face recognition consists in finding the differences in the appearance and constellation of face parts among people. For a dataset with an equal number of faces per person, this method can increase a dataset of  $n$  images to one with  $n^2$  images when using only 2 face parts (we use 5 parts in practice). A dataset like LFW can thus be expanded from a little over 10,000 images to a dataset of 100 million images.

Unlike the existing face synthesis methods [58], [8], [25] which use 3D models, our method is a pure 2D method which is much easier to implement. In addition, our method works on different tasks from [58], [8], [25]. Specifically, the methods [58], [8] are used for facial landmark detection, while ours for face recognition. The approach [25] assumes a relatively large training data (500,000 images) already exists, while we assume the training data (10,000 images) is extremely small.

We experimentally demonstrate that the synthesized large training datasets indeed significantly improve the generalization capacity of CNNs. In our experiments, we generate a training set of 1.5 million images using an initial labeled dataset of only 10,000 images. We improve the face verification rates from 78.97% to 95.77% on LFW using CNNs trained on 10K images and 1.5 million synthetic images respectively. In addition, the proposed face synthesis is also used for NIR-VIS heterogeneous face recognition [27] and improve the rank-1 face identification rate from 17.41% to 85.05%. With the synthetic data, we achieve state-of-the-art performance on both (1) LFW under the “unrestricted, label-free outside data” protocol and (2) CASIA NIR-VIS 2.0 database under rank-1 face identification protocol.

## II. RELATED WORK

Our work relates to three research areas that we briefly review below: face recognition using deep learning methods

(Section II-A), face data collection (Section II-B), and data augmentation and synthesis methods (Section II-C).

### A. Face recognition using deep learning

Here we briefly review recent state-of-the-art face recognition methods based on deep learning.

Since face recognition is a special case of object recognition, good architectures for general object recognition may carry over to face recognition. Schroff et al [35] explored networks that are based on that of Zeiler & Fergus [56] and Inception networks [49]. DeepID3 [43] uses aspects of both Inception networks [48] and the very deep VGG network [40]. Parkhi et al. [29] use the same architecture as the very deep VGG network [40], while Yi et al. [55] use  $3 \times 3$  filters but fewer layers.

DeepFace [50] combines 3D face alignment and CNN for face recognition. Specifically, a 3D model is used for pose normalization, by which all the faces are rotated to the frontal pose. In this way, pose variations are removed from the faces. Then an 8-layer CNN is trained using four million pose-normalized images.

DeepID [46], DeepID2 [42], DeepID2+ [47] all train an ensemble of small CNNs. The input of one small CNN is an image patch cropped around a facial part (face, nose, mouth, etc.). The same idea is also used in [22]. DeepID uses only a classification-based loss to train the CNN, while DeepID2 includes an additional verification-based loss function. To further improve the performance, DeepID2+ adds the loss functions to all the convolutional layers rather than the topmost layer only.

All the above methods train CNNs using large training datasets (500,000 faces or more). To the best of our knowledge, only [13] uses small datasets to train CNNs (only around 10,000 LFW images) and achieves significantly worse performance on the LFW benchmark: 87% vs 97% or higher in [35], [47], [50]. Clearly, sufficiently large training datasets are extremely important for learning deep face representations.

### B. Face dataset collection

Since big data is important for learning a deep face representation, several research groups have collected large datasets with 90,000 up to 2.6 million labeled face images [4], [26], [29], [44], [55]. To achieve this, they collect face images from the internet, by querying for specific websites such as IMDb or general search engines for celebrity names. This data collection process is detailed in [29], [55].

There are, however, two main weaknesses of the existing face data collection methods. First, and most importantly, internet-based collection of large face datasets is limited to visible spectrum images, and is not applicable to collect e.g. infrared face images. Second, the existing collection methods are expensive and time-consuming. It results from the fact that automatically collected face images are noisy, and manual filtering has to be performed to remove incorrectly labeled images [29].

The difficulty of collecting large datasets in some domains, e.g. for infrared imaging, motivates the work presented in

this paper. To address this issue we propose a data synthesis method that we describe in the next section.

### C. Data augmentation and synthesis

The availability of large supervised datasets is the key for machine learning to succeed, and this is true in particular for very powerful deep CNN models with millions of parameters. To alleviate data scarcity in visual recognition tasks, data augmentation has been used to add more examples by applying simple image transformations that do not affect the semantic-level image label, see e.g. [6]. Examples of such transformations are horizontal mirroring, cropping, small rotations, etc. Since it is not always clear in advance which (combinations of) transformations are the most effective to generate examples that improve the learning the most, Paulin et al. [30] proposed to learn which transformations to exploit.

Data augmentation, however, is limited to relatively simple image transformations. Out-of-plane rotations, for example, are hard to accomplish since they would require some degree of 3D scene understanding from a single image. Pose variations of articulated objects are another example of transformations that are non-trivial to obtain, and generally not used in data augmentation methods.

Training models from synthetic data can overcome such difficulties, provided that sufficiently accurate object models are available. Recent examples where visual recognition systems have been trained from synthetic data include the following. Shotton et al. [36] train randomized decision forests for human pose estimation from synthesized 3D depth data. Jaderberg et al. [15] use synthetic data to train CNN models for natural scene text recognition. Su et al. [41] use synthetic images of objects to learn a CNN for viewpoint estimation. Papon and Schoeler [28] train a multi-output CNN that predicts class, pose, and location of objects from realistic cluttered room scenes that are synthesized on the fly. Weinmann et al. [53] synthesize material images under different viewing and lighting conditions based on detailed surface geometry measurements, and use these to train a recognition system using a SIFT-VLAD representation [16]. Ronzantsev et al. [33] use rough 3D models to synthesize new views of real object category instances. They show that this outperforms more basic data augmentation using crops, flips, rotations, etc.

Data synthesis techniques are also used for face analysis. To improve the accuracy of facial landmark detection in the presence of large pose variations [8], [58], 3D morphable face model is used to synthesize face images in arbitrary poses. Similar data synthesis technique is also used for pose-robust face recognition [25]. Unlike 3D solutions, we propose a 2D data synthesis method to solve the problem of training deep CNNs using very limited training data.

## III. SYNTHETIC DATA ENGINE

Human faces are well structured in the sense that they are composed of parts (eyes, nose, mouth, etc.) which are organized in a relatively rigid constellation. Face recognition is conducted implicitly by finding the differences of one or more facial parts and possibly their constellation among people.

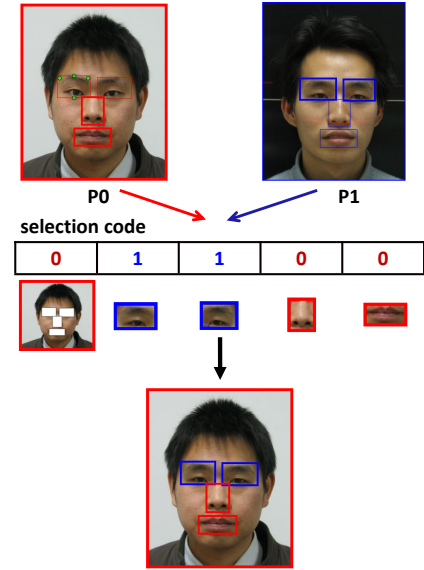


Fig. 1. Schematic illustration of the face synthesis process using five parts: left-eye, right-eye, nose, mouth and the rest. Parent images P0 and P1 (top) are mixed by using the eyes of P1 and the other parts of P0 (middle) to form the synthetic image (bottom).

Motivated by this, our synthetic face images are generated by swapping one or more facial parts among existing “parent” images. In our work we use five face parts: right eye (RE), left eye (LE), nose (N), mouth (M) and the rest (R). See Figure V-A for an illustration. For simplicity, we only consider the synthesis using only two parent images in this work. Our synthesis method can easily be extended, however, to the scenario of more than two parent images.

Suppose that we have an original dataset and let  $\mathcal{S}$  denote the set of subjects in the dataset, and let  $n_i$  denote the number of images of subject  $i \in \mathcal{S}$ . To synthesize an image, we select a tuple  $(i, j, c, s, t)$  where  $i \in \mathcal{S}, j \in \mathcal{S}$  correspond to two subjects that will be mixed, and  $s \in \{1, \dots, n_i\}$  and  $t \in \{1, \dots, n_j\}$  are indices of images of  $i$  and  $j$  that will be used. The bitcode  $c \in \{0, 1\}^5$  defines which parts will be taken from each subject. A zero at a certain position in  $b$  means that the corresponding part will be taken from  $i$ , otherwise it will be taken from  $j$ . There are only  $2^5 - 2 = 30$  valid options for  $b$ , since the codes 00000 and 11111 would correspond to the original images of  $s$  and  $t$  respectively, instead of synthetic ones.

To synthesize a new image, we distinguish the two parent images as the “base” image from which we use the R (the rest) part, and the “injection” image from which one or more parts will be pasted on this base image. Since the size of the facial parts of the two parent images are in general different, we re-size the facial parts of the injection image to that of the base image. The main challenge to implement the proposed synthesis method is to accurately locate the positions of the facial parts. Recently, many efficient and accurate landmark detectors have been proposed. We use four landmarks detected by the method of Zhang et al. [57] to define the rectangular region that corresponds to each face part.

We refer to each choice of  $(i, j, c)$  with  $i \neq j$  as a “virtual subject” which consists of a mix of two existing subjects in the dataset. In total we can generate  $30|\mathcal{S}|(|\mathcal{S}| - 1)/2$  different virtual subjects, and for each of these we can generate  $n_i \times n_j$  samples. Note that if we set  $i = j$  we can in the same manner synthesize  $30n_i(n_i - 1)/2$  new images for an existing subject. We will explore the relative merit of both types of synthetic images in our experiments.

#### IV. FACE RECOGNITION PIPELINE

In this section we describe the different elements of our pipeline for face identification and verification in detail.

##### A. CNN architectures

Face recognition in the wild is a challenging task. As described in Section II-A, the existing deep learning methods highly depend on big training data. Very little research investigates training CNNs using small data. Recently, Hu [13] evaluated CNNs trained using small datasets. Due to the limited training samples, they found the performance of CNNs to be worse than handcrafted features such as high-dimensional features [5] (0.8763 vs 0.9318). In this work, we use a limited training set of around 10,000 images to synthesize a much larger one of around 1.5 million images for CNN training. The synthesized training data captures various deformable facial patterns that are important to improve the generalization capacity of CNNs.

We use two CNN architectures. The first one introduced in [13] has fewer filters and is referred as CNN-S, and the other introduced in [55] is much larger and therefore referred as CNN-L. These two architectures are detailed in Table I. Using the CNN-L model we achieve state-of-the-art performance on the LFW dataset [14] under ‘unrestricted, label-free outside data’ protocol.

##### B. NIR-VIS heterogeneous face recognition

NIR-VIS (near-infrared to visual) face recognition is important in applications where probe images are captured by NIR cameras that use active lighting which is invisible to the human eye [27]. Gallery images are, however, generally only available in the visible spectrum. The existing methods for NIR-VIS face recognition include three steps: (i) illumination pre-processing, (ii) feature extraction, and (iii) metric learning. First, the NIR-VIS illumination differences cause the main difficulty of NIR-VIS face recognition. Therefore, illumination normalization methods are usually used to reduce these differences. Second, to reduce the heterogeneities of NIR and VIS images, illumination-robust features such as LBP are usually extracted. Third, metric learning is widely utilized, aiming at removing the differences of modalities and meanwhile keeping the discriminative information of the extracted features.

In this work, we also follow these three steps that are detailed in Section V-B3. Unlike the existing work that extracts handcrafted features, we learn face representations using two CNN architectures introduced in Section IV-A. To our knowledge, we are the *first* to use deep CNNs for NIR-VIS face

TABLE I  
THE ARCHITECTURES OF THE TWO CNN MODELS WE USED IN OUR EXPERIMENTS.

CNN-L	CNN-S
conv1	
$32 \times 3 \times 3$ , st.1 $64 \times 3 \times 3$ , st.1 x2 maxpool, st.2	$16 \times 3 \times 3$ , st.1 $16 \times 3 \times 3$ , st.1 x2 maxpool, st.2
conv2	
$64 \times 3 \times 3$ , st.1 $128 \times 3 \times 3$ , st.1 x2 maxpool, st.2	$32 \times 3 \times 3$ , st.1 x2 maxpool, st.2
conv3	
$96 \times 3 \times 3$ , st.1 $192 \times 3 \times 3$ , st.1 x2 maxpool, st.2	$48 \times 3 \times 3$ , st.1 x2 maxpool, st.2
conv4	
$128 \times 3 \times 3$ , st.1 $256 \times 3 \times 3$ , st.1 x2 maxpool, st.2	-
conv5	
$160 \times 3 \times 3$ , st.1 $320 \times 3 \times 3$ , st.1 x7 avgpool, st.1	-
fully connected	
Softmax-5000	FC-160 Softmax-5000

recognition. The main difficulty of training CNNs results from the lack of NIR training images which are not available from the Internet. To solve this, we synthesize a big NIR dataset for CNN training.

##### C. Network Fusion

Fusion of multiple networks is a widely used strategy to improve the performance of deep CNN models. For example, in [40], an ensemble of seven networks is used to improve the object recognition performance due to complementarity of the models trained at different scales. Network fusion is also successfully applied to learn face representations. DeepID and its variants [42], [46], [47] train multiple CNNs using image patches extracted from different facial parts.

The heterogeneity of NIR and VIS images is intrinsically caused by the different spectral bands from which they are acquired. The images in both modalities, however, are reflective in nature and affected by illumination variations. Illumination normalization can be used to reduce such variability, at the risk of losing identity-specific characteristics. In this work, we fuse two networks that are trained using the original and illumination-normalized images respectively. This network fusion significantly boosts the recognition rate.

##### D. Metric Learning

The goal of metric learning is to make different classes more separated, and instances in the same class closer. Most approaches learn a Mahalanobis metric

$$d_A^2(x_i, x_j) = (x_i - x_j)^T A (x_i - x_j) \quad (1)$$

which maximizes inter-class discrepancy, while minimizing intra-class discrepancy. Some methods, instead, learn a generalized dot-product of the form

$$d_B^2(x_i, x_j) = x_i^T B x_j \quad (2)$$

Metric learning methods are widely used for face identification and verification. Because identification and verification are two different tasks, different loss functions should be optimized to learn the metric. Joint Bayesian metric learning (JB) [4] and Fisher linear discriminant analysis (LDA) are probably the two most widely used metric learning methods for face verification and identification respectively. In particular, LDA can be seen as a method to learn a metric of the form of Eq. (1), while JB learns a verification function that can be written as a weighted sum of Eq. (1) and (2). In our work we use JB and LDA to improve the performance of face verification and identification respectively.

## V. EXPERIMENTS

### A. Facial data synthesis

Given some face images and their IDs, we define three synthetic strategies: *Inter-Synthesis*, *Intra-Synthesis*, and *Self-Synthesis*. *Inter-Synthesis* synthesizes a new image using two parents from different IDs as shown in Fig. 1. The facial components of an *Intra-Synthesized* face are from different images with the same ID. *Self-synthesis* is a special cause of *Intra-Synthesis*. Specifically, one given image synthesizes new images by swapping facial components of itself and its mirrored images. By virtue of *Self-Synthesis*, one input image can become maximum 32 images which have complementary information. The fusion of features extracted from these 32 images has stronger face representation capacity which is validated in Section V-B2. In the view of NIR-VIS cross-modality, we also define ‘cross-modality synthesis’ which uses images from different modalities to synthesize a new one. Some synthetic images from the CASIA NIR-VIS 2.0 dataset with LSSF [54] illumination normalization are shown in Fig. V-A. The reasons of using LSSF illumination normalization is detailed in Section V-B3. As shown in Fig. V-A, the results of *Intra-Synthesis* method are usually more natural than *Inter-Synthesis* method since the *Intra-Synthesis* method uses the same ID. However, as shown in the right of Fig. V-A, some samples from *Intra-Synthesis* can also be very strange due to large pose variations.

### B. Face recognition

1) *Implementation details*: Our implementation is based on the Caffe open source deep learning toolbox. Before face synthesis, all the raw images are aligned and cropped to size  $100 \times 100$  as in [55] on both datasets. We train our models using images only from LFW and CASIA NIR-VIS2.0 databases. For the CNN-S model on both datasets, we set the learning rate as 0.001, and decrease it by 10 times every 4000 iterations, and stop training after 10K iterations. We practically find dropout is not helpful for the small network, therefore, we train the CNN-S model without dropout. For the CNN-L model on the NIR-VIS dataset, we set the learning rate as 0.01, and decrease it by 10 times every 8000 iterations, and stop training after 20K iterations. For the CNN-L model on the LFW dataset, we set the learning rate as 0.01, and decrease it by 10 times every 120K iterations, and stop training after 200K iterations. We set dropout rate as 0.4 for the pool5 layer

of the CNN-L model. For both CNN-S and CNN-L models, the batch size is 128, momentum is 0.9, and decay is 0.0005. Softmax loss function is used to guide CNN training. The features used in our experiments of CNN-S and CNN-L are FC-160 (160D) and Pool5 (320D), respectively.

#### 2) Face recognition in the wild:

a) *Database and protocol*: Labeled Faces in the Wild (LFW) [14] is a public available dataset for unconstrained face recognition study. It contains 5,749 unique identities and 13,233 face photographs. The training and test sets are pre-defined in [14]. For evaluation, the full dataset is divided into ten splits, and each time nine of them are used for training and the left one for testing. Our work falls in the protocol of “Unrestricted, Label-Free Outside Data” as we use the identity information to train the neural network (softmax loss). Meanwhile, all face images are aligned using a model trained on unlabeled outside data. As a benchmark for comparison, we report the mean and standard deviation of classification accuracy.

b) *Synthetic data generation*: Under LFW protocol, the training set in each fold is different. Therefore, the size of synthetic data and the original raw LFW data in Table II is averaged over 10 folds. We generate 1.5 million training images including 1 million ‘Inter-Syn’ ones and 0.5 million ‘Intra-Syn’ ones. ‘Inter-Syn’ and ‘Intra-Syn’ are defined in Section V-A.

TABLE II  
TRAINING DATA SYNTHESIZED FROM LFW

		IDs	Images	Images/ID
Synthetic	Intra-Syn	5K	500K	100
	Inter-Syn	5K	1M	200
	Total	10K	1.5M	150
Raw		5K	10K	2

c) *Impact of synthetic data*: Table III analyzes the importance of using the synthetic data. First, CNN-S trained using synthetic data (‘Intra-Syn’ and ‘Inter-Syn’) outperforms greatly that trained using raw LFW images, showing the importance of data synthesis. Second, ‘Inter-Syn’ works slightly better than ‘Intra-Syn’ because ‘Inter-Syn’ can capture richer facial variations. Third, combining ‘Inter-Syn’ and ‘Intra-Syn’ works better than either of them because both inter- and intra-personal variations can be captured. Fourth, averaging the features of 32 ‘Self-Syn’ (‘32-Avg’ in Table III and defined in Section V-A) images works consistently better than that of one single test image (‘single’ in Table III). Fifth, CNN-L works consistently better than CNN-S using either raw LFW or synthetic images because deeper architecture has stronger generalization capacity. Last but not least, the metric learning (JB) can further enhance the face recognition performance.

d) *Comparison with the state-of-the-art*: Table IV compares our method with state-of-the-art methods. All methods listed in Table IV except ours use hand-crafted features such as LBP, and this again indicates the hardness of training deep CNNs with small data. In fact, the best deep learning solution [45] recorded in official benchmark achieves 91.75%, and ours is 4% better. In addition, most of state-of-the-art solutions rely on an extremely high dimensional feature vector because they

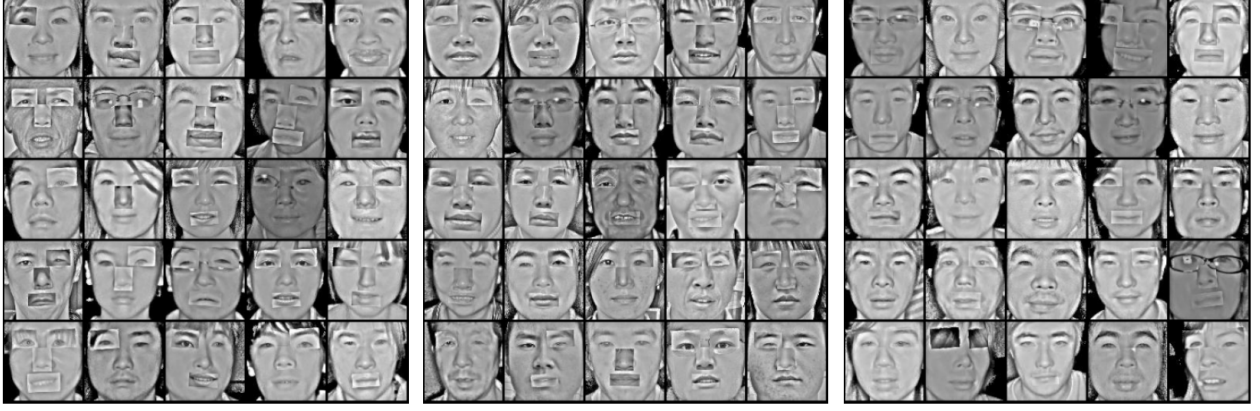


Fig. 2. Left: Inter-Synthesis. Middle: cross-modality Intra-Synthesis. Right: Intra-Synthesis.

TABLE III  
COMPARISON OF SYNTHETIC DATA METHODS ON LFW UNDER ‘UNRESTRICTED, LABEL-FREE OUTSIDE DATA’

Architecture	Metric learning	Training data	single (%)	32-Avg (%)
CNN-S	-	Raw	$78.97 \pm 0.78$	-
		Intra-Syn+Raw	$83.03 \pm 0.56$	$83.93 \pm 0.49$
		Inter-Syn+Raw	$83.18 \pm 0.74$	$84.35 \pm 0.65$
		Intra-Syn+Inter-Syn+Raw	$85.61 \pm 0.71$	$86.98 \pm 0.57$
CNN-L	-	Raw	$85.03 \pm 0.98$	-
	JB [4]	Raw	$87.03 \pm 0.69$	-
	-	Intra-Syn+Inter-Syn+Raw	$94.88 \pm 0.66$	$95.13 \pm 0.53$
	JB [4]	Intra-Syn+Inter-Syn+Raw	$95.32 \pm 0.38$	$95.77 \pm 0.38$

fundamentally employ dense sampling on the face image, in contrast, we just use a 320-dimensional feature vector, which is much more compact than the others.

modalities (NIR and VIS), we applied ‘cross-modality synthesis’ to synthesize new images. The size of synthesized data is detailed in Table V.

TABLE IV  
COMPARISON WITH STATE-OF-THE-ART METHODS ON LFW UNDER ‘UNRESTRICTED, LABEL-FREE OUTSIDE DATA’

Methods	Accuracy (%)
High-dim LBP [5]	$93.18 \pm 1.07$
Fisher vector faces [38]	$93.03 \pm 1.05$
HPEN [59]	$95.25 \pm 0.36$
MDML-DCPs [7]	$95.58 \pm 0.34$
<b>The proposed</b>	<b><math>95.77 \pm 0.38</math></b>

### 3) NIR-VIS face recognition:

a) *Database and protocol:* The largest face database across NIR and VIS spectrum so far is the CASIA NIR-VIS 2.0 face database (CASIA NIR-VIS2.0) [21]. It contains 17,580 images of 725 subjects which exhibit intra-personal variations such as pose and expression. This database includes two views: view 1 for parameter tuning and view 2 including 10 folds for performance evaluation. During test, the gallery and probe images are VIS and NIR images respectively, simulating the scenario of face recognition in the dark environment. The rank 1 identification rate including the mean accuracy and standard deviation of 10 folds are reported.

b) *Synthetic data generation:* We synthesize training samples using the existing images in CASIA NIR-VIS2.0. Because the images of CASIA NIR-VIS2.0 are from two

TABLE V  
TRAINING DATA SYNTHESIZED FROM CASIA NIR-VIS2.0

		IDs	Images	Images/ID
Synthetic	Intra-Syn	357	90K	250
	Inter-Syn	1K	150K	150
	Total	1.4K	240K	170
Raw		357	8.5K	23

### c) Illumination normalization and feature extraction:

Illumination Normalization (IN) methods are usually used to narrow the gap between NIR and VIS images. To investigate the impact of IN, we preprocessed images using three popular IN methods: illumination normalization based on large- and small-scale features (LSSF) [54], DoG filtering-based normalization (DOG) and single-scale retinex (SSR) [17]. We train CNN-S and CNN-L using illumination normalized and non-normalized images. For simplicity, only the images from CASIA NIR-VIS2.0 excluding synthetic ones are used. Fig. 3 shows the face recognition rates at different training iterations using different input images. In Fig. 3 and 4, three IN methods outperform ‘GRAY’ which does not use IN method, showing the effectiveness of IN. Note that LSSF achieves the best performance due to its strong capacity of removing illumination but keep identity information. Same as the LFW results in Section V-B2, CNN-L works better than CNN-S.

TABLE VI  
COMPARISON OF CNN-L VARIANTS AND STATE-OF-THE-ART ON CASIA NIR-VIS2.0 DATABASE

Method		Accuracy (%)	
CNN-L	Training Data	Original	$69.11 \pm 1.21$
		LSSF	$68.97 \pm 1.24$
	Network Fusion	Original+LSSF	$79.96 \pm 1.18$
	Metric Learning	LDA (Original+LSSF)	<b><math>85.05 \pm 0.83</math></b>
State-of-the-art	C-CBFD [24]		$56.6 \pm 2.4$
	C-CBFD+LDA [24]		$81.8 \pm 2.3$
	Dictionary Learning [18]		$78.46 \pm 1.67$

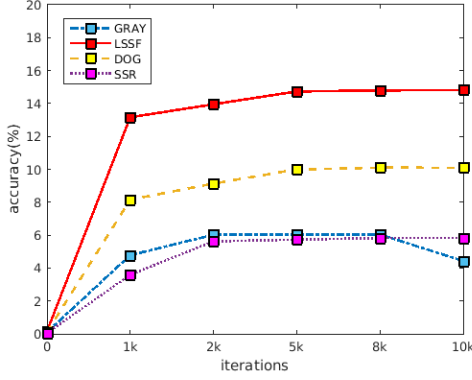


Fig. 3. Comparison of illumination normalization methods using CNN-S

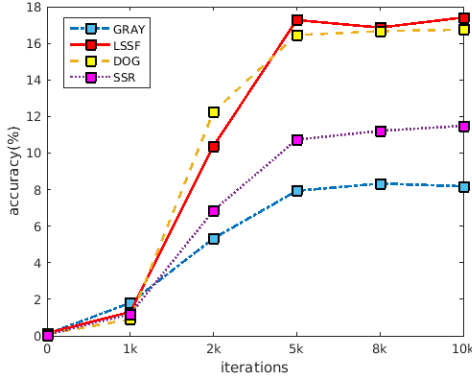


Fig. 4. Comparison of illumination normalization methods using CNN-L

The use of illumination-robust features can effectively narrow the gap between NIR and VIS images. The most commonly used hand-crafted feature is LBP. Since LSSF achieves the best performance, we extract LBP features from LSSF-normalized images and achieve  $12.48 \pm 3.1$  in comparison with  $17.41 \pm 3.76$  by CNN-L learned feature. It shows the superior performance of CNN learned features.

d) *Effects of synthetic data:* To evaluate the effects of synthetic data on a common ground, LSSF is used to preprocess the illumination and CNN-L is applied to learn face representations for all (raw and/or synthetic) input images. In the practice, we find two problems with the synthetic data generated from CASIA NIR-VIS2.0 database: (1) It cannot capture enough facial variations because it only has 357 subjects as shown in Table V. (2) There are much fewer

VIS images than NIR ones. To solve these two problems, we also use the synthetic data generated from LFW images defined in Table II. Table VII compares the results achieved by these two sources of synthetic data. First, the accuracy achieved by using the synthetic data generated from CASIA NIR-VIS2.0 database is  $34.13 \pm 2.13$ , in comparison with  $17.41 \pm 3.76$  without synthetic data. The significant improvement shows the effectiveness of data synthesis. Second, the model trained using raw and synthetic LSSF-normalized LFW images greatly outperforms those synthetic CASIA NIR-VIS2.0 images:  $66.37 \pm 1.45$  vs  $38.45 \pm 2.08$ , although NIR images are completely unseen during training. The reasons are 2-fold: (1) LFW images contains more subjects which can capture more facial variations as analyzed above. (2) LSSF can greatly reduce the gap between NIR and VIS, therefore, LSSF-normalized LFW synthetic images can generalize well to LSSF-normalized NIR images. To further improve the face recognition performance, we trained the network using the synthetic data from both sources (LFW+CASIA NIR-VIS). The face recognition rate is improved from  $66.37 \pm 1.45$  to  $68.97 \pm 1.24$ , showing the importance of bigger synthetic dataset.

TABLE VII  
EVALUATION OF THE IMPACT OF SYNTHETIC DATA ON CASIA NIR-VIS2.0 DATABASE

Training Data			Accuracy(%)
	CASIA NIR-VIS2.0	LFW	
Baseline	Raw	-	$17.41 \pm 3.76$
	Raw+Syn	-	$34.13 \pm 2.13$
Synthetic Data	-	Raw	$38.45 \pm 2.08$
	-	Raw+Syn	$66.37 \pm 1.45$
	Raw+Syn	Raw+Syn	<b><math>68.97 \pm 1.24</math></b>

e) *Comparison with the state-of-the-art:* The CNN-Ls in Table VI are all trained using synthetic LFW data. First, LSSF-normalized and Original LFW synthetic data achieve very comparable performance: 68.97% vs 69.11%. However, the fusion (averaging) of these 2 features can significantly improve the face recognition rates. It shows the fusion can keep the discriminative facial information but remove the illumination effects. Second, not surprisingly, metric learning can further improve the performance. The metric learning method used here is LDA, which is the most widely used one for face identification. Last, Table VI compares the proposed method against the state-of-the-art solutions [24], [18]. [24] uses a designed descriptor that performs better in this dataset compared with other general hand-crafted features, and LDA



can further improve the accuracy. Our method significantly outperforms [24] in the case that LDA is not employed, while it earns 4% advantage with the help of LDA. [18] tries to solve the domain shift between two data sources (NIR and VIS) by a cross-modal metric learning: it assumes that a pair of NIR and VIS images shares the same sparse representation under two jointly learned dictionaries. Our method beats [18] with a 7% margin without such an extra step of dictionary learning.

## VI. CONCLUSION

Recently, convolutional neural networks have attracted a lot of attention in the field of face recognition. However, deep learning methods heavily depend on big training data, which is not always available. To solve this problem in the field of face recognition, we propose a new face synthesis method which swaps the facial components of different face images to generate a new face. With this technique, we achieve state-of-the-art face recognition performance on LFW and CASIA NIR-VIS2.0 face databases. In the future, we will apply this technique to more applications of face analysis.

## ACKNOWLEDGMENT

This work is supported by the ANR project ‘Physionomie’, the European Unions Horizon 2020 research and innovation program under grant No 640891, and Natural Science Foundation of China (No. 61502152). We also gratefully acknowledge NVIDIA Corporation for the donation of the GPUs for this research.

## REFERENCES

- [1] T. Ahonen, A. Hadid, and M. Pietikäinen. Face recognition with local binary patterns. In *Computer vision-eccv 2004*, pages 469–481. Springer, 2004.
- [2] T. Ahonen, E. Rahtu, V. Ojansivu, and J. Heikkilä. Recognition of blurred faces using local phase quantization. In *International Conference on Pattern Recognition*, 2008.
- [3] P. Belhumeur, J. Hespanha, and D. Kriegman. Eigenfaces vs. Fisherfaces: Recognition using class specific linear projection. *PAMI*, 19(7):711–720, 1997.
- [4] D. Chen, X. Cao, L. Wang, F. Wen, and J. Sun. Bayesian face revisited: A joint formulation. In *Computer Vision-ECCV 2012*, pages 566–579. Springer, 2012.
- [5] D. Chen, X. Cao, F. Wen, and J. Sun. Blessing of dimensionality: High-dimensional feature and its efficient compression for face verification. In *Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on*, pages 3025–3032. IEEE, 2013.
- [6] D. Decoste and B. Schölkopf. Training invariant support vector machines. *Machine Learning*, 46:161–190, 2002.
- [7] C. Ding, J. Choi, D. Tao, and L. S. Davis. Multi-directional multi-level dual-cross patterns for robust face recognition. *arXiv preprint arXiv:1401.5311*, 2014.
- [8] Z. Feng, G. Hu, J. Kittler, W. Christmas, and X. Wu. Cascaded collaborative regression for robust facial landmark detection trained using a mixture of synthetic and real images with dynamic weighting. 24(11):3425–3440, 2015.
- [9] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *CVPR*, 2014.
- [10] M. Guillaumin, J. Verbeek, and C. Schmid. Is that you? Metric learning approaches for face identification. In *ICCV*, 2009.
- [11] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. *arXiv preprint arXiv:1512.03385*, 2015.
- [12] X. He, S. Yan, Y. Hu, P. Niyogi, and H.-J. Zhang. Face recognition using Laplacianfaces. *PAMI*, 27(3):328–340, 2005.
- [13] G. Hu, Y. Yang, D. Yi, J. Kittler, W. J. Christmas, S. Z. Li, and T. M. Hospedales. When face recognition meets with deep learning: an evaluation of convolutional neural networks for face recognition. *CoRR*, abs/1504.02351, 2015.
- [14] G. Huang, M. Ramesh, T. Berg, and E. Learned-Miller. Labeled faces in the wild: a database for studying face recognition in unconstrained environments. Technical Report 07-49, University of Massachusetts, Amherst, 2007.
- [15] M. Jaderberg, K. Simonyan, A. Vedaldi, and A. Zisserman. Synthetic data and artificial neural networks for natural scene text recognition. *arXiv preprint arXiv:1406.2227*, 2014.
- [16] H. Jégou, M. Douze, C. Schmid, and P. Pérez. Aggregating local descriptors into a compact image representation. In *CVPR*, 2010.
- [17] D. J. Jobson, Z.-U. Rahman, and G. A. Woodell. Properties and performance of a center/surround retinex. *IEEE Trans. Image Processing*, 6(3):451–462, 1997.
- [18] F. Juefei-Xu, D. Pal, and M. Savvides. Nir-vis heterogeneous face recognition via cross-spectral joint dictionary learning and reconstruction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 141–150, 2015.
- [19] S. Kong, J. Heo, B. Abidi, J. Paik, and M. Abidi. Recent advances in visual and infrared face recognition – a review. *CVIU*, 97(1):103 – 135, 2005.
- [20] A. Krizhevsky, I. Sutskever, and G. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012.
- [21] S. Z. Li, D. Yi, Z. Lei, and S. Liao. The casia nir-vis 2.0 face database. In *Computer Vision and Pattern Recognition Workshops (CVPRW), 2013 IEEE Conference on*, pages 348–353. IEEE, 2013.
- [22] J. Liu, Y. Deng, and C. Huang. Targeting ultimate accuracy: Face recognition via deep embedding. *arXiv preprint arXiv:1506.07310*, 2015.
- [23] D. Lowe. Distinctive image features from scale-invariant keypoints. *IJCV*, 60(2):91–110, 2004.
- [24] J. Lu, V. E. Liong, X. Zhou, and J. Zhou. Learning compact binary face descriptor for face recognition. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 2015.
- [25] I. Masi, A. T. Tran, J. T. Leksut, T. Hassner, and G. Medioni. Do we really need to collect millions of faces for effective face recognition? *arXiv preprint arXiv:1603.07057*, 2016.
- [26] D. Miller, E. Brossard, S. M. Seitz, and I. Kemelmacher-Shlizerman. Megaface: A million faces for recognition at scale. *arXiv preprint arXiv:1505.02108*, 2015.
- [27] S. Ouyang, T. Hospedales, Y.-Z. Song, and X. Li. A survey on heterogeneous face recognition: Sketch, infra-red, 3d and low-resolution. *arXiv preprint arXiv:1409.5114*, 2016.
- [28] J. Papon and M. Schoeler. Semantic pose using deep networks trained on synthetic rgb-d. *arXiv preprint arXiv:1508.00835*, 2015.
- [29] O. Parkhi, A. Vedaldi, and A. Zisserman. Deep face recognition. In *BMVC*, 2015.
- [30] M. Paulin, J. Revaud, Z. Harchaoui, F. Perronnin, and C. Schmid. Transformation pursuit for image classification. In *CVPR*, 2014.
- [31] P. Pinheiro and R. Collobert. From image-level to pixel-level labeling with convolutional networks. In *CVPR*, 2015.
- [32] A. Rozantsev, V. Lepetit, and P. Fua. On rendering synthetic images for training an object detector. *Computer Vision and Image Understanding*, 2015.
- [33] A. Rozantsev, V. Lepetit, and P. Fua. On rendering synthetic images for training an object detector. *CVIU*, 137:24 – 37, 2015.
- [34] J. Sánchez, F. Perronnin, T. Mensink, and J. Verbeek. Image classification with the Fisher vector: Theory and practice. *IJCV*, 105(3):222–245, 2013.
- [35] F. Schroff, D. Kalenichenko, and J. Philbin. Facenet: A unified embedding for face recognition and clustering. In *CVPR*, 2015.
- [36] J. Shotton, T. Sharp, A. Kipman, A. Fitzgibbon, M. Finocchio, A. Blake, M. Cook, and R. Moore. Real-time human pose recognition in parts from single depth images. *Communications of the ACM*, 56(1):116–124, 2013.
- [37] T. Sim, S. Baker, and M. Bsat. The CMU Pose, Illumination, and Expression (PIE) database. In *IEEE International Conference on Automatic Face and Gesture Recognition*, 2002.
- [38] K. Simonyan, O. M. Parkhi, A. Vedaldi, and A. Zisserman. Fisher Vector Faces in the Wild. In *British Machine Vision Conference*, 2013.
- [39] K. Simonyan and A. Zisserman. Two-stream convolutional networks for action recognition in videos. In *NIPS*, 2014.
- [40] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015.
- [41] H. Su, C. Qi, Y. Li, and L. Guibas. Render for cnn: Viewpoint estimation in images using cnns trained with rendered 3d model views. In *ICCV*, 2015.
- [42] Y. Sun, Y. Chen, X. Wang, and X. Tang. Deep learning face representation by joint identification-verification. In *Advances in Neural Information Processing Systems*, pages 1988–1996, 2014.
- [43] Y. Sun, D. Liang, X. Wang, and X. Tang. Deepid3: Face recognition with very deep neural networks. *arXiv preprint arXiv:1502.00873*, 2015.
- [44] Y. Sun, X. Wang, and X. Tang. Hybrid deep learning for face



- verification. In *ICCV*, 2013.
- [45] Y. Sun, X. Wang, and X. Tang. Hybrid deep learning for face verification. In *Computer Vision (ICCV), 2013 IEEE International Conference on*, pages 1489–1496. IEEE, 2013.
  - [46] Y. Sun, X. Wang, and X. Tang. Deep learning face representation from predicting 10,000 classes. In *Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on*, pages 1891–1898. IEEE, 2014.
  - [47] Y. Sun, X. Wang, and X. Tang. Deeply learned face representations are sparse, selective, and robust. *arXiv preprint arXiv:1412.1265*, 2014.
  - [48] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. *arXiv preprint arXiv:1409.4842*, 2014.
  - [49] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *CVPR*, 2015.
  - [50] Y. Taigman, M. Yang, M. Ranzato, and L. Wolf. DeepFace: Closing the gap to human-level performance in face verification. In *CVPR*, 2014.
  - [51] M. A. Turk and A. P. Pentland. Face recognition using eigenfaces. In *Computer Vision and Pattern Recognition*, pages 586–591, 1991.
  - [52] K. Weinberger and L. Saul. Distance metric learning for large margin nearest neighbor classification. *JMLR*, 10:207–244, 2009.
  - [53] M. Weinmann, J. Gall, and R. Klein. Material classification based on training data synthesized using a BTF database. In *ECCV*, 2014.
  - [54] X. Xie, W.-S. Zheng, J. Lai, P. C. Yuen, and C. Y. Suen. Normalization of face illumination based on large-and small-scale features. *IEEE Trans. Image Processing*, 20(7):1807–1821, 2011.
  - [55] D. Yi, Z. Lei, S. Liao, and S. Z. Li. Learning face representation from scratch. *arXiv preprint arXiv:1411.7923*, 2014.
  - [56] M. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. In *ECCV*, 2014.
  - [57] Z. Zhang, P. Luo, C. Loy, and X. Tang. Facial landmark detection by deep multi-task learning. In *ECCV*, 2014.
  - [58] X. Zhu, Z. Lei, X. Liu, H. Shi, and S. Z. Li. Face alignment across large poses: A 3d solution. *arXiv preprint arXiv:1511.07212*, 2015.
  - [59] X. Zhu, Z. Lei, J. Yan, D. Yi, and S. Z. Li. High-fidelity pose and expression normalization for face recognition in the wild. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 787–796, 2015.



**Yongxin Yang** is a PhD student at the department of Electronic Engineering and Computer Science, Queen Mary University of London (QMUL). His research interests include machine learning (Transfer Learning, Domain Adaptation, Multi-Task Learning and Deep Learning) and computer vision.



**Timothy Hospedales** is a Senior Lecturer in Computer Science at QMUL. He is head of the Applied Machine Learning Lab at QMUL, and a member of the Risk and Information Management (RIM) group. His research focuses on probabilistic modelling and machine learning, particularly life-long, transfer and active learning.



**Guosheng Hu** is a postdoctoral researcher in THOTH team, INRIA Grenoble Rhone-Alpes, France. He received his PhD degree in the Centre for Vision, Speech and Signal Processing, University of Surrey, UK in 2015. His research interests include pattern recognition, biometrics, machine learning and graphics. His PhD topic is ‘face analysis using 3D morphable models. He is an IEEE member.



**Xiaojiang Peng** is a postdoctoral researcher in THOTH team, INRIA Grenoble Rhone-Alpes, France. He received his PhD degree in Computer Science at Southwest Jiaotong University in 2014. His research focus is in the areas of action recognition and deep learning.



**Jakob Verbeek** received a PhD degree in computer science in 2004 from the University of Amsterdam, The Netherlands. After being a postdoctoral researcher at the University of Amsterdam and at INRIA Rhone-Alpes, he has been a full-time researcher at THOTH group, INRIA, Grenoble, France, since 2007. His research interests include machine learning and computer vision, with special interest in applications of statistical models in computer vision.